



The Lucene Scalpel

Tuning Low-Level Segments for Multilingual Concurrency

Divyanshu Mishra

AI & DevOps Enthusiast

Agenda

01

Setting the Stage

What Lucene segments are & why non-engineers should care



02

The Multi-Lingual Challenge

How language diversity breaks your search performance



03

The Tuning Toolkit

Practical levers: merge policy, codecs, analyzers



04

Real-World Results

Before/after benchmarks from production systems



05

Your Action Plan

Configs, checklist & next steps you can use Monday morning



Who Is This Talk For?

No PhD required. If you search, store, or decide about search — this is for you.



Backend / Search Engineer

Actionable configs & code you can apply this week



Solutions Architect

Mental models for explaining trade-offs to stakeholders



Engineering Manager / Lead

Why segment tuning is a lever, not just maintenance



Product / Platform Owner

Business impact of latency at multi-lingual scale

01 · Setting the Stage — What Is a Lucene Segment?

Think of it like this:

A library with thousands of books. OpenSearch is the library. A shard is a floor. A segment is a bookshelf — self-contained, sorted, never re-sorted after shelving.



Immutable

Once written, a segment never changes.
Updates create new segments.



Self-contained

Each has its own inverted index, metadata & term data.



Merged over time

Background merging keeps your library tidy & fast.

⚡ An OpenSearch shard is just a collection of segments — tuning them is tuning your entire search experience.

02 · The Multi-Lingual Challenge

The core problem:

Different languages don't just need different words — they need fundamentally different indexing strategies. One size fits none.



Analyzer Wars

Arabic stemmer + ICU tokenizer + English analyzer compete for the same flush thread. Winner: slowness.



Term Explosion

Arabic/Hebrew morphology produces 3–5× more index terms per document than English. Bigger segments, slower merges.



CJK vs Latin Gap

Chinese/Japanese use character n-grams → thousands of tiny segments.
Latin uses word tokens → fewer, larger. Mixed = chaos.



Concurrency Bottleneck

Every new segment batch forces readers to open a fresh IndexSearcher.
More languages = more batches = more overhead.

02 · A Story You Might Recognize

Scenario: E-commerce platform, 12 languages, 50M products

English queries: P95 = 45ms ✓ Arabic queries: P95 = 340ms ✗ Mixed-language: P95 = 580ms ✗ ✗

What was actually happening under the hood:

1

Index time

Arabic products created 4× more segments than English equivalents during bulk ingest

2

Merge pressure

TieredMergePolicy (default settings) treated all segments equally — CJK segments kept fragmenting

3

Query time

Each Arabic/mixed query opened 80–120 segment readers vs 12–15 for English-only queries

4

The fix

Per-language merge policy + field-level analyzers reduced Arabic P95 from 340ms → 68ms

MERGE

1. TieredMergePolicy

Control how many segments merge at once — tune per language index tier

ANALYZE

2. Per-Field Analyzers

Assign the right analyzer to each language field — never share one globally

CODEC

3. Codec Selection

best_compression for CJK/RTL fields; default for Latin to balance speed

MEMORY

4. RAM Buffer Sizing

Larger buffer = fewer, bigger flush segments = fewer merge ops overall

OPS

5. Force Merge Strategy

Run forceMerge after bulk ingest, never during live search traffic

LATENCY


6. NRT Refresh Interval

Hot tier: 1s refresh. Warm: 10s. Cold: 60s. Don't use one value for all.

03 · How to Decide What to Tune First

Use this 2-question diagnostic before touching any config:

Q1: What's your language mix?

 **Mostly one language**

Default settings OK — focus on shard sizing

 **2–3 languages**

Per-field analyzers + gentle Tiered Merge tuning

 **4+ languages / CJK+RTL**

Full per-language index strategy needed

Q2: Where is your pain?

 **Slow queries**

Reduce segment count → merge policy & forceMerge

 **Slow indexing**

Increase RAM buffer, async translog, codec tweak

 **Both**

Start with merge policy — it affects both paths

04 · Real-World Results — Before vs After

62%

**P99 Query Latency
Reduction**

Multi-lingual mixed query

3.1×

**Indexing Throughput
Improvement**

Arabic + English bulk ingest

41%

**Segment Count
Reduction**

After TieredMerge tuning

Query Latency by Language (ms, P95) — Same Cluster, Same Hardware



opensearch_index.json

```
{
  "settings": {
    "index": {
      "codec": "best_compression",
      "merge.policy": {
        "max_merge_at_once": 5,
        "segments_per_tier": 8,
        "expunge_deletes_allowed": 10
      },
      "refresh_interval": "5s",
      "translog.durability": "async"
    }
  },
  "mappings": {
    "properties": {
      "content_ar": {
        "analyzer": "arabic" },
      "content_zh": {
        "analyzer": "smartcn" },
      "content_en": {
        "analyzer": "english" }
    }
  }
}
```

Pre-Flight Checklist

- Know your language mix before choosing merge policy
- Use per-field analyzers — never one global analyzer
- Enable best_compression for CJK or RTL-heavy indexes
- Monitor segment count via `_cat/segments` API weekly
- Schedule `forceMerge` during off-peak hours only
- Tune `refresh_interval` per tier (hot: 1s, warm: 10s)
- Always test P99 latency — not just averages
- Benchmark each language separately first

Key Takeaways — What to Remember

1

Segments are the unit of search performance

Understanding them unlocks every other tuning decision — for engineers and architects alike

2

Language shapes your index topology

CJK and RTL languages create fundamentally different segment pressure — treat them differently

3

One analyzer for all is a trap

Per-field analyzers with dedicated merge profiles unlock 40–60% latency gains in production

4

Measure P99, not averages

Your worst-case user is a multi-lingual query during merge pressure — that's who you're tuning for

Thank You!

The Lucene Scalpel: Tuning Low-Level Segments for Multi-Lingual Concurrency

[Learn More](#)

OpenSearch Docs → opensearch.org/docs/latest/search-plugins/

 **OpenSearchCon**
EUROPE

